



TITLE:

An analysis of the feasibility and benefits of GPU/multicore acceleration of the Weather Research and Forecasting model

AUTHOR(S):

Vanderbauwhede, Wim; Takemi, Tetsuya

CITATION:

Vanderbauwhede, Wim ...[et al]. An analysis of the feasibility and benefits of GPU/multicore acceleration of the Weather Research and Forecasting model. *Concurrency and Computation: Practice and Experience* 2016, 28(7): 2052-2072

ISSUE DATE:

2016-05

URL:

<http://hdl.handle.net/2433/218036>

RIGHT:

This is the accepted version of the following article: [Vanderbauwhede, W., and Takemi, T. (2016) An analysis of the feasibility and benefits of GPU/multicore acceleration of the Weather Research and Forecasting model. *Concurrency Computat.: Pract. Exper.*, 28: 2052–2072], which has been published in final form at <http://doi.org/10.1002/cpe.3522>. This article may be used for non-commercial purposes in accordance with Wiley Terms and Conditions for Self-Archiving.; The full-text file will be made open to the public on 8 APR 2017 in accordance with publisher's 'Terms and Conditions for Self-Archiving'; この論文は出版社版ではありません。引用の際には出版社版をご確認ください。; This is not the published version. Please cite only the published version.

An Analysis of the Feasibility and Benefits of GPU/Multicore Acceleration of the Weather Research and Forecasting Model

Wim Vanderbauwhede^{1*}, Tetsuya Takemi²

¹*School of Computing Science, University of Glasgow, Glasgow, UK*

²*Disaster Prevention Research Institute, Kyoto University, Uji, Kyoto, Japan*

SUMMARY

There is a growing need for ever more accurate climate and weather simulations to be delivered in shorter timescales, in particular to guard against severe weather events such as hurricanes and heavy rainfall. Due to climate change, the severity and frequency of such events – and thus the economic impact – are set to rise dramatically. Hardware Acceleration using GPUs or FPGAs could potentially result in much reduced run times or higher accuracy simulations.

In this paper, we present the results of a study of the Weather Research and Forecasting (WRF) model undertaken in order to assess if GPU and multicore acceleration of this type of Numerical Weather Prediction (NWP) code is both feasible and worthwhile. The focus of this paper is on acceleration of code running on a single compute node through offloading of parts of the code to an accelerator such as a GPU. Therefore, in this work we deployed WRF on a single node consisting of a multicore CPU and GPGPU. We did not use a cluster of nodes in our experiment as the scalability and performance of WRF in a cluster is determined by the MPI subsystem, and our GPU/multicore acceleration approach is entirely orthogonal to this. In other words from a cluster perspective an accelerated node is simply a faster node, and the MPI performance and scalability of WRF has already been investigated in detail on a variety of systems [1, 2, 3]. The governing equation set of the WRF model is based on the compressible, non-hydrostatic atmospheric motion with multi-physics processes. We put this work into context by discussing its more general applicability to multi-physics fluid dynamics codes: in many fluid dynamics codes the numerical schemes of the advection terms are based on finite differences between neighboring cells, similar to the WRF code. For fluid systems including multi-physics processes, there are many calls to these advection routines. This class of numerical codes will benefit from hardware acceleration.

We studied the performance of the original code of the WRF model and created a simple model for comparing multicore CPU and GPU performance. Based on the results of extensive profiling of representative WRF runs, we focused on the acceleration of the scalar advection module. We discuss the implementation of this module as a data-parallel kernel in both OpenCL and OpenMP.

We show that our data-parallel kernel version of the scalar advection module runs up to $7\times$ faster on the GPU compared to the original code on the CPU. However, as the data transfer cost between GPU and CPU is very high (as shown by our analysis), there is only a small speed-up ($2\times$) for the fully integrated code.

In order to carry out this research, we also developed an extensible software system for integrating OpenCL code into large Fortran code bases such as WRF, which is one of the main contributions of our work. We discuss the system to show how it allows to replace sections of the original codebase with their OpenCL counterparts with minimal changes – literally only a few lines – to the original code.

Our final assessment is that even with the current system architectures, accelerating WRF – and hence also other, similar types of multi-physics fluid dynamics codes – with a factor of up to five times is definitely and achievable goal.

Copyright © 2014 John Wiley & Sons, Ltd. *Concurrency Computat.: Pract. Exper.* (2014)
Prepared using *cpeauth.cls* DOI: 10.1002/cpe

Accelerating multi-physics fluid dynamics codes including NWP codes is vital for its application to weather

CONCURRENCY AND COMPUTATION: PRACTICE AND EXPERIENCE

Concurrency Computat.: Pract. Exper. 2014; **00**:3–35

Received ...

Published online in Wiley InterScience (www.interscience.wiley.com). DOI: 10.1002/cpe

KEY WORDS: General-Purpose computation on Graphics Processing Units (GPGPU), Parallelization of Simulation, Large Scale Scientific Computing

1. BACKGROUND

There is a growing need for ever more accurate climate and weather simulations to be delivered in shorter timescales, in particular to guard against severe weather events such as hurricanes and heavy rainfall. Due to climate change, the severity and frequency of such events – and thus the economic impact – are set to rise dramatically[4, 5]. Hardware Acceleration using GPUs or FPGAs could potentially result in much reduced run times or higher accuracy simulations. As climate change will result in more, and more severe extreme weather events, faster, more accurate predictions of extreme weather events are needed [6, 7]. Understanding climate change itself requires growing amounts of computational power [8].

1.1. The Weather Research and Forecasting Model

The Weather Research and Forecasting Model[†] (WRF) [9, 10, 11] is a state-of-the-art mesoscale numerical weather prediction system (NWP) intended both for forecasting and atmospheric research. It is an Open Source project, created by a partnership of the US National Oceanic and Atmospheric Administration (NOAA), the National Center for Atmospheric Research (NCAR), and more than 150 other organizations and universities; it is used by a large fraction of weather and climate scientists worldwide. The WRF code base is written in Fortran-90 and is both complex and extensive (about a million lines of code). The governing equation set of the WRF model is based on the compressible, non-hydrostatic atmospheric motion with multiple physics processes such as cloud and precipitation, boundary-layer turbulence, land-ocean-air interaction, radiative transfer in the atmosphere, and energy transfer at the surface. The finite difference method is used to discretize

*Correspondence to: School of Computing Science, University of Glasgow, G12 8QQ Glasgow, UK

[†]<http://www.wrf-model.org>

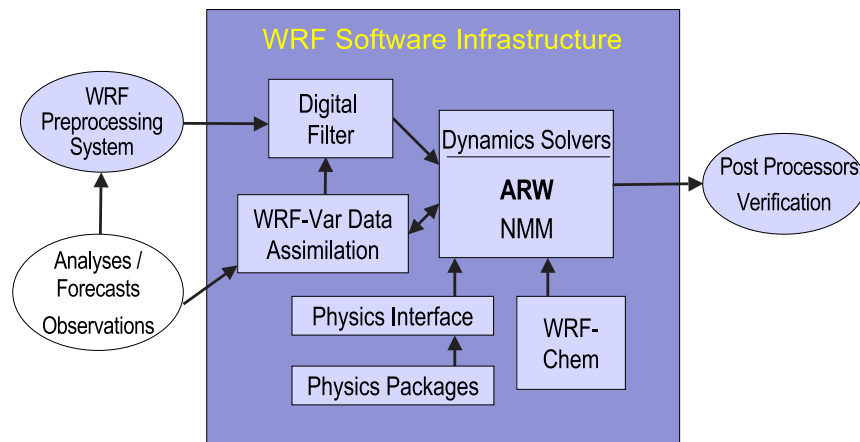


Figure 1. WRF-ARW system components (from [10])

the governing equations of the WRF model. These discretized equations are integrated in time to obtain time-dependent atmospheric motion and physical states. Owing to the multiple physical processes that determine the atmospheric motion field, the number of the prognostic variables of the WRF model is quite large compared to a simple computational fluid dynamics (CFD) model that consists of the Navier-Stokes equation and the mass continuity equation. The large number of the prognostic variables in the three dimensions is a severe computational constraint, which requires a high-performance computational resource.

In this paper we focus on the advanced research version of WRF, called WRF-ARW (Advanced Research WRF) [10] which features very high resolution and is being used to explore ways of improving the accuracy of simulation of severe weather events, e.g. tropical cyclones such as hurricanes and typhoons, tornadoes, windstorms, and heavy rainfall events. The WRF-ARW system components are depicted in Figure 1. The most computationally intensive components are the Dynamics Solvers and the Physics Packages.

1.2. Previous Work on GPU Acceleration of WRF

Previous work in GPU acceleration of WRF is discussed [12]. This work dealt with the WRF Single-Moment 5-class (WSM5) microphysics kernel[‡] (one of the Physics packages). It was an experimental, stand-alone implementation and is not included in the standard WRF distribution. However, it demonstrated the potential for accelerating weather physics codes on GPUs. Furthermore, both implementations are in CUDA. As CUDA is a proprietary technology, we prefer to use OpenCL instead. OpenCL has the advantage that it can be deployed on GPUs and multicore CPUs of different versions. Recently OpenCL support for FPGAs has become available, this is a very promising technology for NWP. Several other NWP codes have been adapted for GPU [13, 14, 15]. However, because of its size and complexity, a full GPU port of WRF has not yet been undertaken.

1.3. OpenCL Programming

OpenCL [16] was developed by the Khronos Group in 2008 as an open standard for parallel programming of heterogeneous systems and is finding increasing adoption amongst providers of multicore CPUs and GPUs (e.g. AMD, Intel, ARM) and FPGAs (Altera). It provides an API for control and data transfer between the host and device (typically the host CPU and a GPU) and a language for kernel development. Contrary to proprietary solutions such as Nvidia's CUDA and Microsoft's DirectX, OpenCL is open and cross-platform, so that it can be deployed on different operating systems (Linux, OS X, Windows) and hardware architectures (multicore CPUs, GPUs, FPGAs). The OpenCL API is defined for C and a C++. In practice, the API is quite fine grained and verbose and requires a lot of boiler plate code to be written. Consequently, it is not straightforward to integrate OpenCL in existing codes, especially for non-computing scientists.

[‡]<http://www.mmm.ucar.edu/wrf/WG2/GPU/WSM5.htm>

1.4. The OclWrapper Library

To facilitate the integration of the OpenCL code into the existing code base, we developed the OclWrapper library[§] which supports C, C++ and Fortran-95. The library wraps the OpenCL platform, context and command queue into a single object, with a much smaller number of calls required to run an OpenCL computation. As it is a thin wrapper, the additional abstraction comes at no cost in terms of features: the OpenCL API is completely accessible.

The OclWrapper library consists of several components:

The OclWrapper C++ Class This class abstracts the OpenCL concepts of Platform, Context, Device and Command Queue using a single object. Using C++ features such as templates, polymorphic functions and default arguments, it provides a greatly simplified interface that is suitable for the majority of OpenCL applications. However, as the object includes all the lower-level OpenCL objects, all low-level OpenCL features are still accessible without overhead.

The oclWrapper Fortran-95 Library This library provides the *oclWrapper* Fortran module, which offers a subroutine-based interface to the C++ OclWrapper class. As Fortran does not offer polymorphic subroutines, the library provides individual functions for manipulating multi-dimensional arrays of various types.

The oclBuilder SCons Library To build the OclWrapper, we use the SCons[¶] build system, a replacement for Make that allows to write very complex build scripts in Python. The oclBuilder library makes it possible to write a build script for OclWrapper in a few lines.

The use of the library is illustrated below on a simple C++ OpenCL example.

```
// Create wrapper for default device and single kernel  
OclWrapper ocl(srcfilename, kernelname, opts);  
  
// Create read and write buffers
```

[§]<https://github.com/wimvanderbauwhede/OpenCLIntegration>

[¶]<http://scons.org/>

AN ANALYSIS OF THE FEASIBILITY AND BENEFITS OF GPU/MULTICORE ACCELERATION OF WRF 7

```
cl::Buffer rbuf = ocl.makeReadBuffer(sz);  
cl::Buffer wbuf = ocl.makeWriteBuffer(sz);  
  
// Transfer input data to device  
ocl.writeBuffer(rbuf, sz, warray);  
  
// Set up index space  
ocl.enqueueNDRange(globalrange, localrange);  
  
// Run kernel  
ocl.runKernel(wbuf, rbuf).wait();  
  
// Read output data from device  
ocl.readBuffer(wbuf, sz, rarray);
```

Without the wrapper, the same program would be about a hundred lines of code, and each individual call would have many arguments. Furthermore, the use of wrapper in Fortran is equally straightforward:

```
use oclWrapper  
integer(8) :: rbuf, wbuf  
real :: dimension (ims:ime, kms:kme, jms:jme) rarray  
real :: dimension (ims:ime, kms:kme, jms:jme) warray  
integer :: globalrange, localrange  
  
! Create wrapper for default device and single kernel  
call oclInit(srcfilename, kernelname)  
  
! Create read and write buffers  
call oclMake3DFloatArrayReadBuffer(rbuf, sz, rarray)  
call oclMake3DFloatArrayWriteBuffer(wbuf, sz)  
  
! Transfer input data to device  
call oclWrite3DFloatArrayBuffer(rbuf, sz, rarray)  
  
! Run kernel over index space  
oclRun(globalrange, localrange)  
  
! Read output data from device  
oclRead3DFloatArrayBuffer(wbuf, sz, warray)
```

Moreover, because the Fortran `oclWrapper` is a module, it stores the `OclWrapper` object globally, so that the subroutine calls can be issues in different program units.

1.5. Hardware Performance Indicators

1.5.1. Computational Performance Indicator We used several different systems for this work. The host system used for the main experiments was based on an Intel Xeon E5-2640 CPU (dual-processor, 6 cores/chip, 2 threads/core). This processor has 256-bit AVX SIMD, so a smart compiler will do up to 8 floating point operations in parallel. The GPU was an Nvidia GeForce GX480. It has 15 Compute Units with 32 Processing Elements each.

Two other systems, a 48-core AMD Opteron 6176-SE system (four 12-core processors) and a Tesla C2070 GPU hosted on a 24-core AMD Opteron 6174 system (two 12-core processors), were used for additional experiments (see Table I for full details). We define the (single precisions floating point) computational performance indicator as

$$CPI = \#threads \times SIMD\ width \times clock\ freq$$

This figure is directly proportional to flops, but more easy to obtain. We define “threads” as the product of the number of cores/compute units and their hyperthreading capability, and “vector size” as either the SIMD vector size or the number of processing elements per compute unit. The CPIs for our platforms are shown in Table I.

From the table we see that purely in terms of computation, under optimal circumstances, the GeForce GPU can be at best $1.4\times$ faster than the Intel CPU; the CPI of the Tesla GPU is only 3% higher than that of the Intel CPU. If the memory bandwidth is the limiting factor, the achievable speed-up for the application running on the GPU would be $4.2\times$. The total achievable speed-up is limited by the data transfer rate between host memory and GPU memory, and the overhead for control of the GPU. According to our measurements we achieve about 2 GB/s reading from the GPU and 8 GB/s writing to the GPU. In Section 5 we present the detailed discussion of the cost of data transfer and computation.

	#cores	vector size	Clock speed (GHz)	CPI	Memory BW (GB/s)
CPU: Intel Xeon E5-2640	24	8	2.5	480	42.6
GPU: Nvidia GeForce GX480	15	32	1.4	672	177.4
CPU: AMD Opteron 6176 SE	48	4	2.3	441.6	42.7
CPU: AMD Opteron 6174	24	4	2.1	201.6	42.7
GPU: Nvidia Tesla C2070	14	32	1.1	492.8	144

Table I. Specifications of hardware platforms used in this work

1.5.2. Communication Bandwidth Limit on Achievable Performance In Section 5 we present the detailed discussion of the cost of data transfer and computation. In general, we can analyse the achievable speed-up as a function of the computational speed up (which in its term depends on the CPI and the memory bandwidth) and the data transfer speed. Figure 2 shows a generic graph which can be used to assess the performance of an algorithm.

What the graph shows is the achievable speed up as a function of the CPU compute time relative to the data transfer time, with the GPU/CPU computational speed-up as a parameter. For example, if the computation on the CPU takes 100ms, and the data transfer 1000ms, then there can be no speed-up, no matter how fast the GPU computes. On the other hand, if the CPU takes 1000ms and the transfer time is 100ms, then with a GPU/CPU computational speed-up of $5\times$ the total speed-up = $1000 / (100+1000/5) = 3.3\times$. In the legend, "pipelined" means that the computations and data transfers overlap in pipelined fashion, which can improve performance when processing a stream of data, and if the transfer and compute times are of the same order.

2. METHODOLOGY

To assess the feasibility of GPU acceleration of WRF, we did the following:

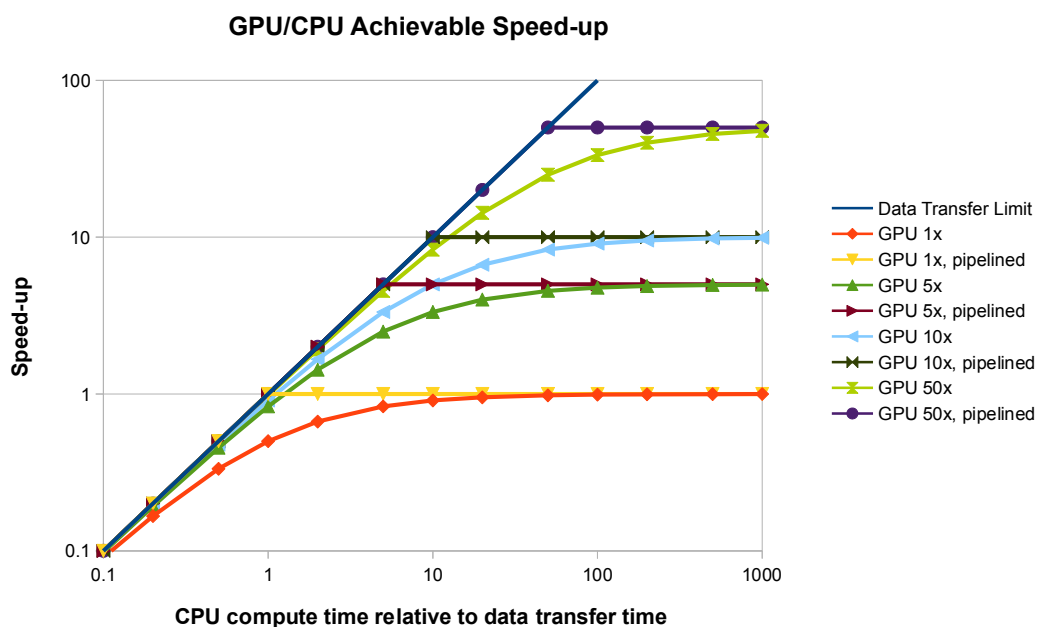


Figure 2. Achievable speed-up from offloading work to the GPU

1. Performance evaluation of the current WRF software using MPI and OpenMP
2. Profiling of WRF runs with a number of different configurations
3. Selection of code portions suitable for acceleration via data-parallel computation
4. Implementation of the code in OpenCL
5. Performance evaluation of the OpenCL kernel
6. Integration of the OpenCL kernel into the WRF code

3. WRF-ARW PERFORMANCE ANALYSIS

3.1. Settings of the WRF simulation

The version of the WRF-ARW model used here is version 3.4, which was released in April 2012. The fifth-order upwind biased scheme is used for the discretization of the advection terms in the horizontal direction, the third-order upwind biased scheme is used for the discretization of the

AN ANALYSIS OF THE FEASIBILITY AND BENEFITS OF GPU/MULTICORE ACCELERATION OF WRF 11

advection terms in the vertical direction, and the third-order Runge-Kutta scheme is used for the time integration of the governing equations [10].

The case investigated in the present numerical simulations is the severe tornado case that occurred in Tsukuba, Japan, a suburban area north of the Tokyo metropolitan region, on 6 May 2012. This tornado spawned severe damages in Tsukuba and its surroundings and was rated as the F3 scale by the Fujita tornado damage scale. The simulation settings are set up with the full physics modules implemented in WRF. Since the meteorological case chosen here is a tornado that was generated by a well-developed cumulonimbus cloud system, one of the most important physical processes is a cloud and precipitation process, so called a microphysics process. The WRF Single-Moment 6-class (WSM6) microphysics module [17] is used for the microphysics parametrization, because this scheme is one of the most sophisticated single-moment schemes and is successful in dealing with convective storms in moist regions such as East Asia [17].

The WRF model has a capability of setting multiple computational domains nested in a larger domain. The present study explores the computational performance of two cases of domain settings: one is a single domain, and the other case sets triple nested domains. For the single domain case the horizontal grid spacing is 5 km, while for the triple domain case the grid spacings are 2.5 km, 500 m, and 100 m.

3.2. MPI versus OpenMP

To establish the baseline performance, we carried out a number of experiments of WRF runs with MPI and OpenMP. The first sets of results (Figure 3) was obtained on the AMD 24-core system, using the GNU Fortran compiler (*gfortran*). It compares OpenMP with MPI performance for a small domain size without nesting:

```
e_we = 100,  
e_sn = 100,  
  
e_vert = 27,
```

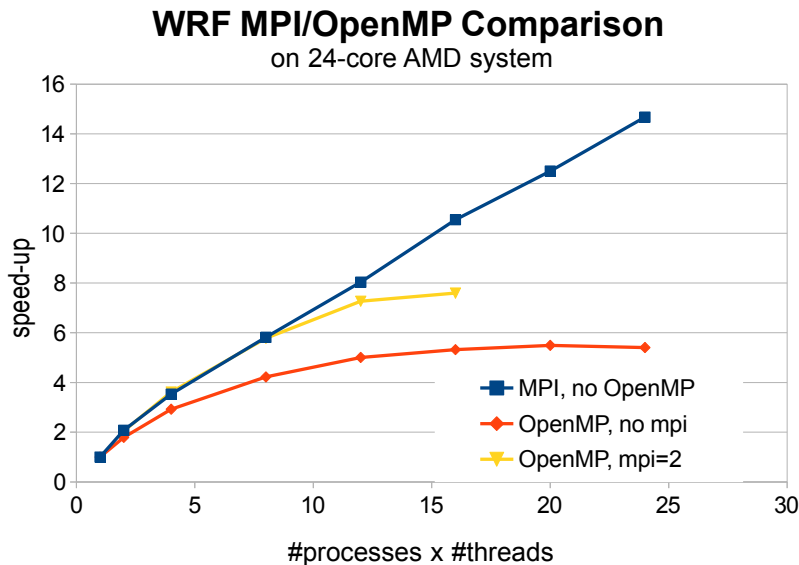



Figure 3. WRF performance with MPI and OpenMP, domain size 100x100x27, on the 24-core AMD system

where e_{we} , w_{sn} , and e_{vert} are the sizes of computational grid in the east-west, the north-south, and the vertical direction. The second set of results is obtained on the 48-core AMD system, using the Intel Fortran compiler (*ifort*). A third set of results was obtained on a 12-core (24-thread) Intel Xeon system, it shows OpenMP performance for varying numbers of threads. Both the second and third experiment simulate a larger domain with nesting:

```
e_we = 500, 301, 501,
e_sn = 500, 301, 501,
e_vert = 60, 60, 60
```

The other WRF settings are identical. The second and third set are shown in Figure 4.

The figures show the speed-up as a function of the number of *parallel processes* × *threads*. Missing points indicate that the simulation failed to complete.

The conclusions from these experiments are clear:

- First, performance varies considerably across platforms.
- Second, for large domains, the speed-up saturates at about half the number of physical threads.
- Finally, MPI outperforms or matches OpenMP for all cases.

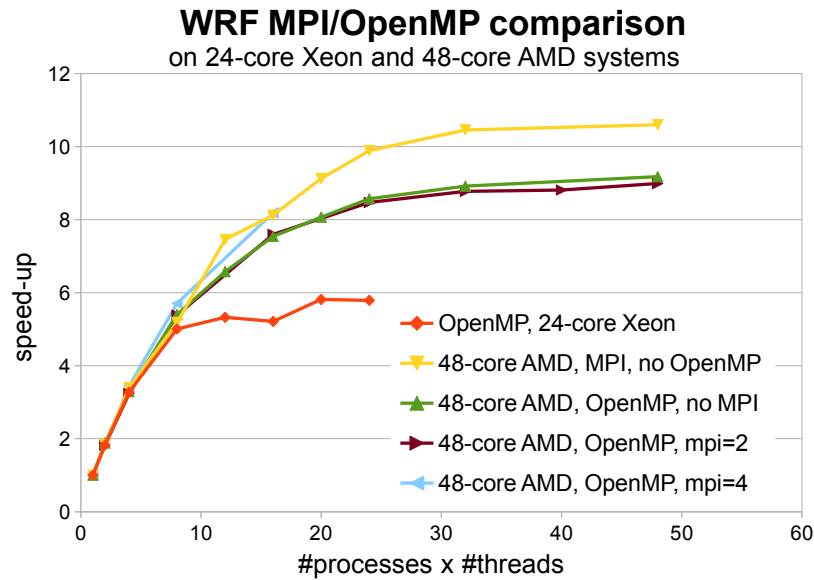


Figure 4. WRF performance with MPI and OpenMP, domain size 500x500x60 with nesting, on the 24-core Intel and 48-core AMD systems

These findings are in line with other studies, e.g. [18]. The poor OpenMP performance is due to the sub-optimal use of OpenMP in WRF: the use of many shared variables results in frequent locking. This is a result of the decision not to rewrite the code for OpenMP, but to rely only on insertion of pragmas.

The observed behavior constitutes a problem for GPU acceleration, and would need to be addressed: in current systems, the GPU can only be accessed by a single process at a time. As MPI creates separate processes, access to the GPU would be serialized. Furthermore, each computation on the GPU would be on the portion of the total memory space used by the MPI process, rather than on the full memory space. So either all processes would have to copy their memory space to the process controlling the GPU, or the GPU would have to be called sequentially by each process in turn. Either way, as a result the overhead of accessing the GPU would dominate the performance, and the net result would be a slow-down rather than a speed-up. With effective use of OpenMP, it should be possible to match or even better the MPI performance (as intrinsically OpenMP has a lower overhead). The host code would be a single process with a single memory space and could easily

functionality	modules	test1	test2
dynamics		71	65
	advection: advect_scalar advect_scalar_pd	28	20
	small-step	17	17
	big-step	17	18
	other	9	10
physics		14	21
	microphysics	7	12
	other	7	9

Table II. Contributions of various parts of WRF to total run time (%)

and effectively interact with the GPU. An alternative solution would be to rewrite the MPI code to include an additional process which would have access to the combined memory space and control the GPU. As the OS uses copy-on-write, this approach should also result in good performance, and might even be preferred.

3.3. WRF Run Profiling

We profiled the two typical WRF runs (on a 256x256x32 domain) using the Shark sampling profiler (on OS X 10.6.8). The conclusion of these experiments was that most of the time is spent in the dynamic core and, to a lesser extent, the physics modules (See Table II)

Together, dynamics and physics constitutes about 85% of the total run time. This time is divided across a large number of calls to different routines, so there is no “quick win”. However, the contributions of advection, big-step and small-step routines and microphysics already account for 70% of the total run time, so these parts of the model constitute a logical focus for acceleration. These findings are in line with those for the COSMO model [14]. Note that accelerating 70% of

AN ANALYSIS OF THE FEASIBILITY AND BENEFITS OF GPU/MULTICORE ACCELERATION OF WRF 15

the code with $5\times$ results in a speed-up of $2.3\times$, a speed-up of $10\times$ on 85% of the code would result in $4.3\times$. Furthermore, the structure of the dynamics kernels is similar in terms of the required approach to parallelisation, so that by studying one of the kernels we can infer the behaviour of the other kernels.

4. OPENCL KERNEL FOR SCALAR ADVECTION

A relatively large part of the run time for WRF is spent in the scalar advection routines *advect_scalar* and *advect_scalar_pd*. These routines are part of the dynamic core (*dyn_em*), and no previous GPU implementations have been reported. As all the advection routines are all similar in structure, the OpenCL version of *advect_scalar* can serve as a template for the other routines.

4.1. Approach

The original WRF kernel for scalar advection consists of a number of nested loops over i, j, k , where typically the inner loops are guarded by *if*-statements. Also, the code uses arrays to store all intermediate results. First, we translated the code to C, using *F2C_ACC* [13]. We then analyzed the conditionals and replaced all run-time *if*-statements that are actually run-time constants with preprocessor *if*-statements. This is important for GPU kernels as run-time branching of a thread in a single warp will lead to stalling of the threads that do not execute the selected branch. The next step was an analysis of the loop structures and boundaries, resulting eventually in a single, unified nested loop with conditionals inside.

This approach is not appropriate for single-threaded code as the new code executes more statements because the conditions are evaluated for every combined loop iteration. However, for data parallel execution, the placement of the conditionals as in the original code would not result in reduced run times, only in thread stalling.

We then replaced the intermediate arrays with local variables and removed some loop dependencies by computing “ahead of time”. Finally, we merged the loops into a single loop, and then used this loop’s range as the index space (the global *NDRange*). The local *NDRange* was set to

the k -range. We experimented with different approaches and values for global and local ranges, but found that the above configuration was optimal.

In terms of effort, the total elapsed time to parallelise the kernel, port it to OpenCL and validate it was about one month, for an experienced computing scientist. The total project including the software engineering required to integrate the OpenCL code seamlessly into the original WRF code took two months.

4.2. Implementation

The structure of the kernel is shown in Algorithm 1. The array *ranges_boundaries_degrade* contains the various ranges (*ims*, *ime* etc.) and computed boundaries (*i_start*, *i_end* etc.) and conditions (*degrade_xs*,...) for the computation. It is more efficient to pass these to the kernel as an array than as individual arguments. The *zero_tendency* argument is used to determine the part of the kernel to be executed (first zero the *tendency* array, then compute the new values). Every thread only uses a small portion of the *field* array (typically ± 3 grid points in every dimension), therefore we copy the required values to a local arrays for *i,j* and *k*. The functions *calc_tendency_** contain the advection computations for the x, y and z dimensions.

4.3. Verification

In order to verify that our OpenCL kernel code produces the same results as the original Fortran code, we employed a testing approach where both codes are run in succession on identical input values, and the results computed by each are compared at run time using a set of comparison functions. To account for differences in rounding errors arising from the different order in which the floating point instructions are executed, we allowed an error of $5 \cdot 10^{-6}$ (2 bits).

4.4. Performance Evaluation

The rewritten code is intended for data-parallel execution on a GPU. It is therefore expected that the code will run slower than the original sequential code. This is confirmed by our measurements (Figure 5): the data-parallel kernel, when run sequentially, is about $4\times$ slower than the original

Algorithm 1 Structure of the OpenCL scalar advection kernel

```
__kernel void advect_scalar (  
    __global float *tendency,  
    __global const float *field,  
    // ... other data ...  
    __constant const int *ranges_boundaries_degrade,  
    __constant const int *zero_tendency  
) {  
    int gl_id = get_global_id(0);  
    if (zero_tendency[0] != 0) {  
        // zero the tendency array  
        tendency[gl_id] = 0.0;  
    } else {  
        // assign ranges_boundaries_degrade to local variables for convenience  
        // calculate the ranges for i, j, k  
        // calculate i,j,k from the global index  
        // create a local copy lfield of the field entries for i,j,k needed for calculating the fluxes  
        // read tendencies for x, y and z  
        float tend_ikj = 0.0;  
        // calculate the tendencies for x, y and z  
        if (j >= j_start_y && j <= j_end_y) {  
            tend_ikj = calc_tendency_y_l(...);  
        }  
        if (i >= i_start_x && i <= i_end_x) {  
            tend_ikj = calc_tendency_x_l(...);  
        }  
        tend_ikj = calc_tendency_z_l(...);  
        // write the result to main memory  
        tendency[...] = tend_ikj;  
    }  
}
```

Fortran code compiled with *gfortran* or the equivalent C code compiled with *gcc*, with optimization *-O3*. This is expected because, as detailed in Section 4.1, the rewritten code executes many more instructions than the original code, as a result of moving conditional branches from the outside to the inside of loops. However, when compiled with the PGI Fortran compiler *pgfortran*, with optimizations *-fast -fastsse -Mipa=fast*, the original Fortran code is $5\times$ faster compared to the *gfortran* binary. This is a result of the better vectorization performance of the GPI compiler: the Xeon E5-2640 CPU has 256-bit AVX vectors, so it can in principle handle 8 single-precision floating-point operations in parallel. The *gcc* compiler does not vectorize the code, hence the observed performance difference. We also tested the effect of the auto-parallelization option *-Mconcur* with various sub-options, but this does not result in performance improvement.

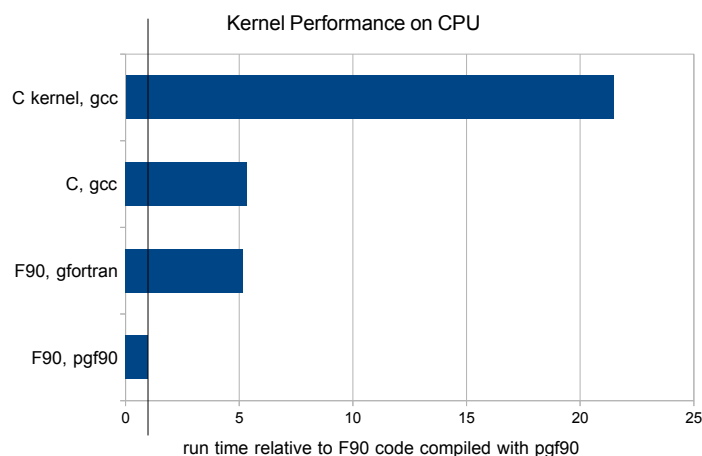


Figure 5. Kernel performance compared to original code, both running single-threaded on CPU.

Figure 6 shows the performance of the actual OpenCL kernel (which is essentially the same code as the C kernel but parallelized using the OpenCL framework), relative to the performance of *gcc/gfortran*, which we chose as the reference because it is available on all CPU platforms we used. We observe a speed-up of about $12\times$ on the GeForce GPU (the reasons for the lower performance on the Tesla GPU are discussed in Section 5).

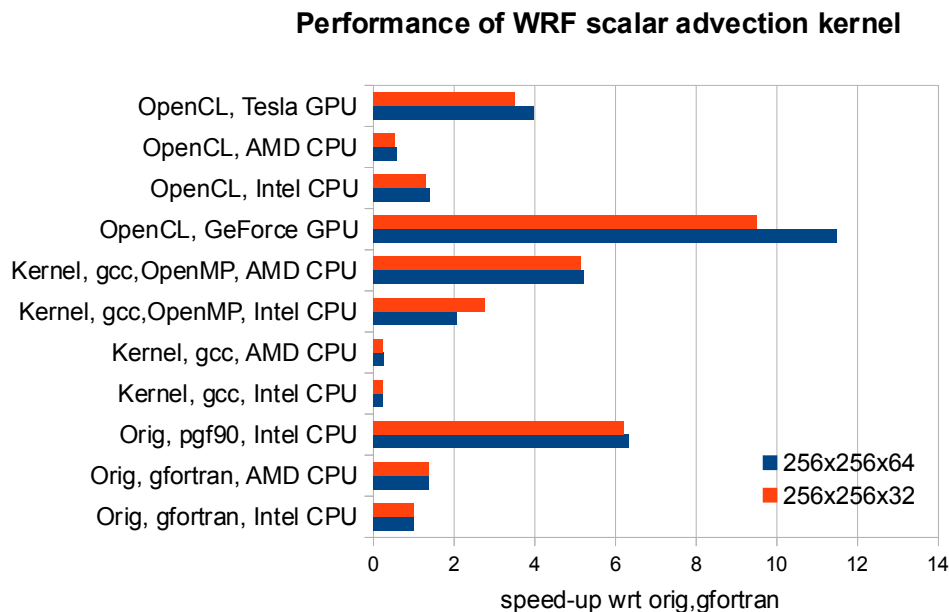


Figure 6. OpenCL kernel performance

We also evaluated the performance when parallelizing the kernel execution with OpenMP (Figure 7). On the AMD system, the speed-up for 24 threads was $21\times$. This The performance on the Intel CPU was also very good but saturated at $8\times$ at 12 threads, this shows that the performance of a hyperthreaded core is less good than that of two separate cores for this type of code. This is because all threads are busy most of the time: hyperthreading works essentially by allowing more than one thread to run per core, but this mechanism is only effective when the threads are stalling a lot of the time: under such circumstances, without hyperthreading the CPU would idle, with hyperthreading the CPU is used by one thread while the other is stalled.

4.5. Discussion of Kernel Performance

At first sight it might seem from these results that the GPU acceleration is hardly worth it: the OpenCL code deployed on the GeForce GTX 480 GPU is only about twice as fast as the original code when compiled with the PGI compiler. However, it is important to realize that accelerating

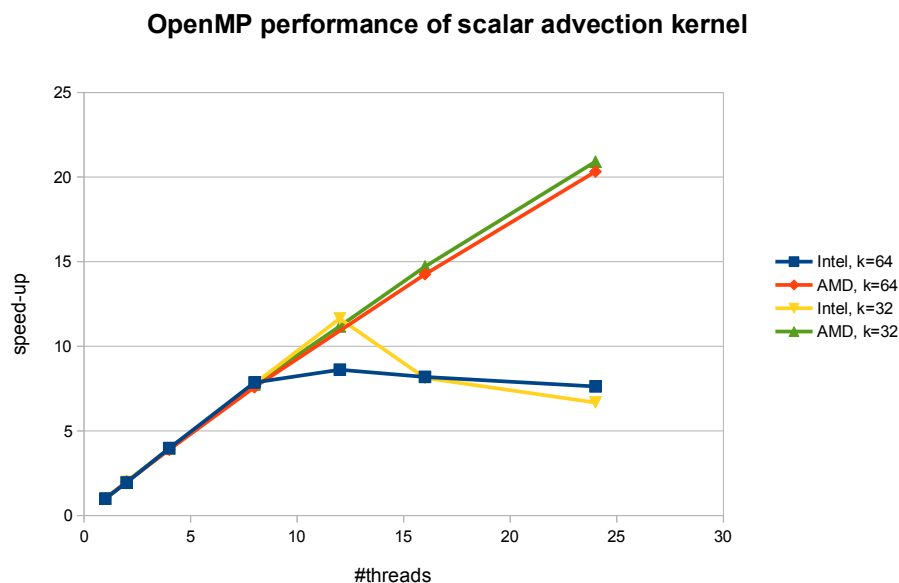


Figure 7. OpenMP performance of the scalar advection kernel, $(i,j,k)=(256,256,32)$ and $(256,256,64)$

only the scalar advection kernel would not speed up WRF execution anyway, as it accounts for only about 10% of the run time.

As explained in Section 3.3, a large portion of the code base must be accelerated to the GPU to achieve considerable speed-ups of the total application (Amdahl's law). As we will see from the analysis in Section 5, under those circumstances GPU acceleration can result in considerable performance increase.

5. GPU RUN TIME ANALYSIS

The specifications in Section 1.5 provide a good guideline for the achievable performance; however, to get a clear picture, in this Section we present an analysis of the performance of the *advect_scalar* kernel on a GeForce GTX480 and a Tesla C2070 GPU.

AN ANALYSIS OF THE FEASIBILITY AND BENEFITS OF GPU/MULTICORE ACCELERATION OF WRF21

5.1. Experiments

We performed the following experiment on the `advect_scalar` kernel, with a domain size of $256 \times 256 \times 64$:

- For each run of the GPU, the host:
 - writes 4 buffers of the domain size (16MB) to the GPU memory, and a number of smaller buffers, total transferring about 64.25 MB
 - calls the GPU twice: first to zero the tendency array, then to compute the new tendencies
 - reads the new tendency array, 16MB
- The host performed 100 runs in a loop and recorded the aggregate run time.
- This experiment was repeated 20 times

The run time contributions for both GPUs are shown in Figure 8. In this figure, “compute only” means no data transfer from host memory to GPU memory; “data transfer only” means that the kernel is called but performs no computation; “no zeroing” means that the call to zero the tendency array is skipped, so the GPU is called only once per run.

We also investigated the influence of the data size on the performance. The experiment was the same as above, but we varied the domain size as follows:

```
i, j: 32, 64, 128, 256, 512, 1024  
k: 32, 64
```

Figure 9 shows the speed-up of the OpenCL GPU code compared to the original code with the GNU and PGI fortran compilers.

There are several interesting points about these results:

5.1.1. Influence of Host System on Transfer Time The first is the difference in transfer time: the AMD/Tesla system takes almost $4 \times$ longer than the Intel/GeForce system: the transfer bandwidth for the Intel/GeForce system is 2.8GB/s (100 transfers of 64MB in 2.3s), which is reasonably close to the top performance of 4GB/s for a 16-lane PCI express with 2.5GT/s transfer rate; however, the

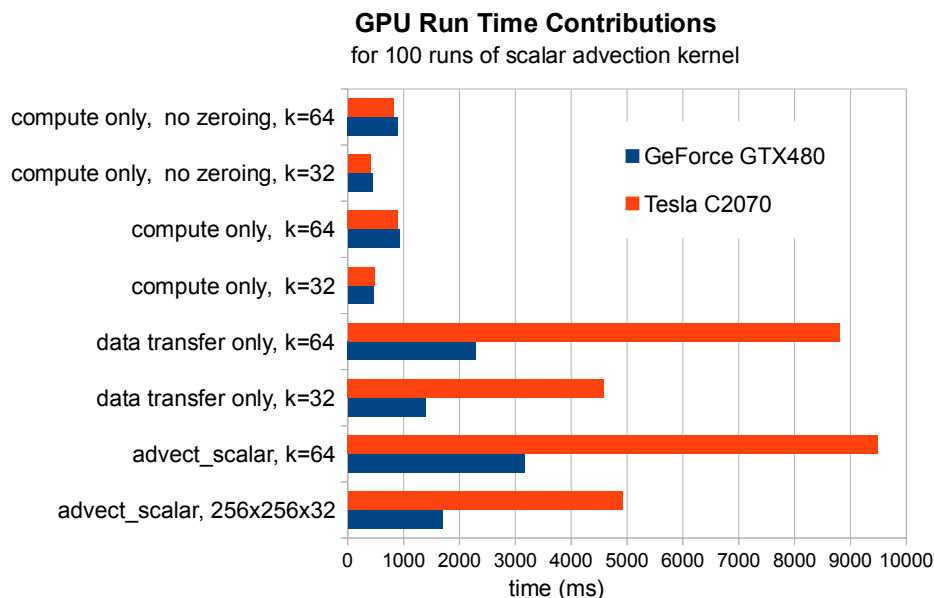


Figure 8. Run time contributions for GeForce GTX 480 GPU on Intel host and Tesla C2070 on AMD host

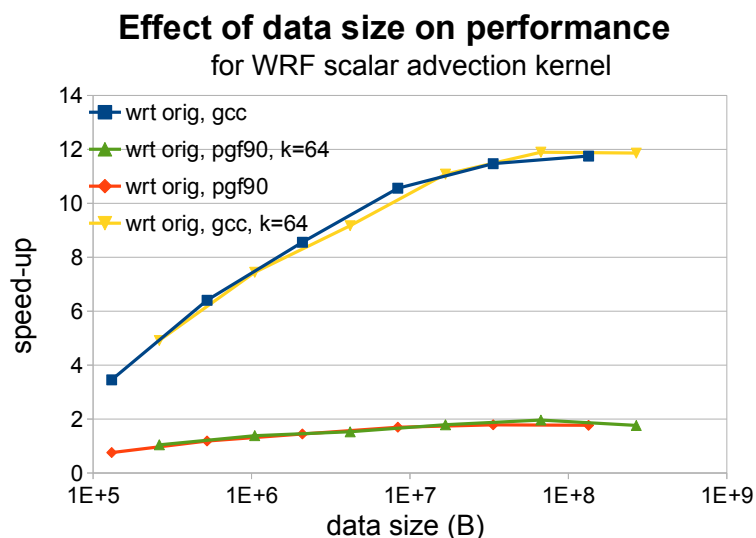


Figure 9. Influence of data size on performance

AMD/Tesla system only reaches 730MB/s. Looking closer at the the PCIe specs of both systems, the only difference is in the latency: both systems have a 16-lane PCI Express v2, 2.5GT/s, but the more recent Intel system has a latency $< 256ns$, whereas the AMD system has a latency $< 1\mu s$. The AMD system also has a significantly lower memory bandwidth (see Figure 10). Another factor that

AN ANALYSIS OF THE FEASIBILITY AND BENEFITS OF GPU/MULTICORE ACCELERATION OF WRF23

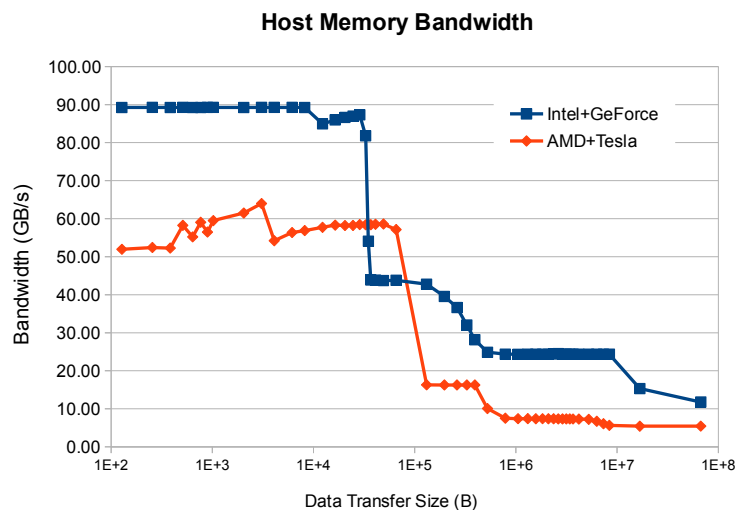


Figure 10. Host memory bandwidth

most likely influences the transfer time is the smaller cache size; in any case, for a transfer size of $> 1MB$ the AMD CPU's memory bandwidth is $2.5\times$ lower than the Intel CPU's.

It should be noted that a more modern system with PCIe v3 should be capable of 10GT/s, so the achievable performance of the GPU computation would be considerably better: for the scalar advection kernel, the total run time would be reduced by a factor of two.

5.1.2. GPU Compute Performance The second observation is that the Tesla GPU computes the kernel about as fast as the GeForce: this is not really surprising when comparing the specs of both GPU cards: the main difference is in the amount of on-board memory. NVidia mentions that the floating point performance of the Tesla cards is better than that of the “consumer cards”, but that applies only to double-precision floating point. As the WRF uses single precision, the much cheaper GeForce card is the better choice.

5.1.3. Influence of the Data Size We see from Figure 9 that the speed-up of the GPU increases for larger data sizes. The reasons for this behavior are twofold: on the one hand, the fixed cost for starting the GPU is relatively less important for larger data transfers. On the other hand, the computation

grows more than linear with data size, so the cost if the data transfer is less dominant for larger data sizes. There are other factors, e.g. cache mis-alignment is less important on larger transfers. All these factors contribute to the observed behavior.

If we compare only the compute performance of the GPU with the original code on the Intel CPU (compiled with the PGI compiler), we see that the GPU is up to $7\times$ faster (Figure 11).

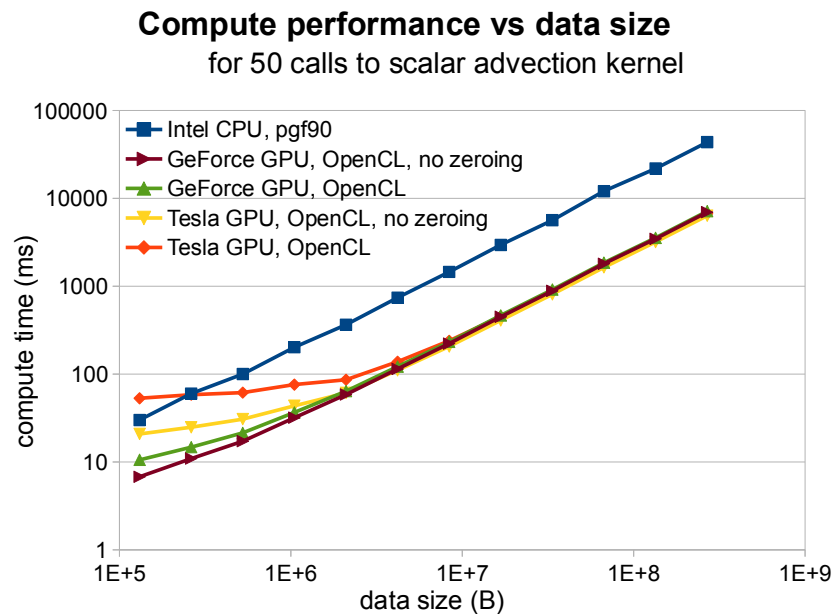


Figure 11. Compute performance comparison of original code on Intel CPU, compiled with the PGI Fortran compiler, to OpenCL code on Tesla and GeForce GPUs.

In summary, the conclusions of this analysis are:

- The host-GPU link is the main bottleneck, and care must be taken in the choice of the host platform, in particular memory bandwidth and PCIe latency.
- The GPU needs to work on large data sizes for optimal performance. For our kernel, performance is optimal for data sizes $> 64MB$.
- If the complete computation was performed on the GPU, this would yield a speed-up of $5\times$ – $10\times$ compared to the original code compiled with the PGI compiler.

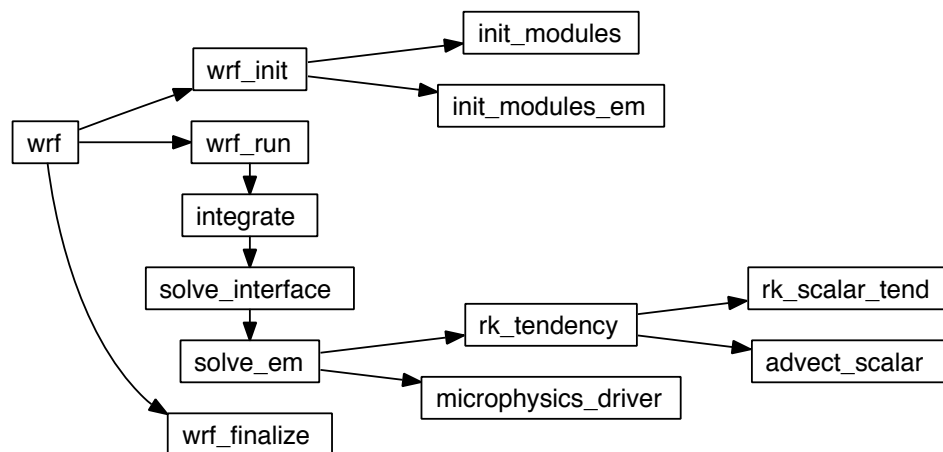


Figure 12. WRF code structure

6. INTEGRATION OF THE OPENCL CODE INTO WRF

As part of this research we developed a strategy for integrating OpenCL code into large Fortran codebases such as WRF. We created a simplified Fortran OpenCL wrapper library^{||} to facilitate the integration.

6.1. Overview of WRF Code Structure

The WRF code structure for the dynamic kernel, omitting details, is shown in Figure 12. To integrate the OpenCL host code code, we added a single *use* statement and a single call to *wrf_init* in *main/module_wrf_top.F*:

```

subroutine wrf_init( no_initl )
! ... other use statements ...
use module_init_ocl
! ... original wr_init code ...
call advect_scalar_init_ocl_grid ( head_grid )

```

^{||}<https://github.com/wimvanderbauwhede/OpenCLIntegration>

```
end subroutine wrf_init
```

This is the only change to the initialisation code.

The change to the dynamic kernel code is also minimal: in the subroutine *rk_tendency* (in *dyn_em/module_em.F*) we added:

```
use module_advect_scalar_ocl
```

And we replace the call to *advect_scalar* using the preprocessor by a call to *advect_scalar_ocl*:

```
call advect_scalar_ocl
```

We could have given the OpenCL routine the same signature as the original *advect_scalar* routine but for clarity we prefer to have a separate name.

The actual new source code is entirely contained in two new modules, *init_ocl* and *advect_scalar_ocl*.

The first module contains two subroutines, *advect_scalar_init_ocl_grid* and *advect_scalar_init_ocl*. The first routine (Algorithm 2) is called in *wrf_init* as shown above. It is essentially a wrapper routine which extracts information from the *head_grid* datastructure and passes it on to *advect_scalar_init_ocl* (Algorithm 3), which performs the OpenCL framework initialisation. The main actions in this routine are loading and compiling the kernel, creating the buffers and setting the kernel arguments.

The second module (Algorithm 4) runs the OpenCL scalar advection kernel on the GPU. Its main actions are writing the data to the GPU, running the GPU and reading back the data. Note that the GPU is run twice, once to zero the tendencies and once to compute the new tendencies.

To extend this work, rather than making separate calls to all the different dynamics and physics routines, the aim is to create an OpenCL version of the full *solve_em* subroutine. Then we can simply replace *solve_em* by *solve_em_ocl* in quite the same way as above.

Algorithm 2 OpenCL initialisation wrapper to extract info from grid

```
subroutine advect_scalar_init_ocl_grid ( grid )  
  
    use module_domain  
  
    ! Variable declarations  
  
    ...  
  
    call get_ijk_from_grid ( ... )  
  
    ! Loop range computations  
  
    ...  
  
    call nl_get_time_step ( 1, time_step )  
  
    ! Call actual OpenCL initialisation routine  
  
    call advect_scalar_init_ocl ( ... )  
  
end subroutine
```

7. DISCUSSION

The main research questions we set out to answer in this work was: is hardware acceleration of the Weather Research and Forecasting model on GPUs feasible and worthwhile?

First, we studied the code and performed experiments on the current parallel performance using MPI and OpenMP. We found that the WRF OpenMP performance is sub-optimal as it performs worse than MPI, whereas in principle OpenMP should have considerably smaller overhead. We analyzed the reasons for this behavior, and concluded that to amend it is a major effort. We also observed that the MPI behavior is strongly sub-linear and saturates typically when 50% of the available hardware threads have been used, and at a performance of less than half the maximally achievable performance. In other words, it is in principle possible to speed up WRF considerably.

Then we profiled WRF runs to identify the most important routines in terms of run time. Our findings, confirmed by other authors, are that the dynamics and physics account for the majority of the WRF run time. As there has been previous work on acceleration of physics modules, we focused on the dynamics, and in particular we chose the scalar advection module as the target for our study as it is the dominant routine in terms of run time. By implementing the GPU kernel and evaluating its performance, we get more detailed answers to the questions of feasibility and pay-off.

Algorithm 3 OpenCL initialisation routine

```
subroutine advect_scalar_init_ocl (...)
    use oclWrapper

    ! Load the kernel source and compile for the platform

    srcstr='advect_scalar_ocl.cc'
    kstr='advect_scalar'
    call oclInit(srcstr,kstr)

    ! Set up the ranges

    oclGlobalRange = gl_range
    oclLocalRange = 0 ! NullRange

    ! Create the buffers

    call oclMakeWriteBuffer(tendency_buf,jkmfsz)
    call oclMakeFloatArrayReadBuffer(field_buf,jik_sz,field)
    !...

    ! Set the Kernel arguments

    call oclSetFloatArrayArg(0, tendency_buf )
    !...

    ! Assign to module array for convenience

    oclBuffers(1) = tendency_buf
    !...

end subroutine ! advect_scalar_init_ocl
```

Before discussing the parallelization of the WRF scalar advection kernel, we want to discuss the capabilities of the hardware platforms used in this work. There are a lot of unrealistic expectations considering the achievable performance of multicore CPUs and GPUs. To help understand the performance of these systems we defined indicators for the compute capability and memory bandwidth. From these indicators, we concluded that for computation-dominated code, the theoretical speed-up achievable by running the code on the GPU is quite small: moving the code from the Intel Xeon to the Nvidia Geforce could result in a speed-up of $1.4\times$; moving the code from the AMD host CPU to the Nvidia Tesla C2070 can at best provide a 3% speed-up. As noted above,

Algorithm 4 OpenCL driver code for scalar advection kernel

```

subroutine advect_scalar_ocl (...)

    use oclWrapper

    !...

! Write buffers to GPU memory

    call oclWriteBuffer(field_buf, jik_sz, field);

    !...

! First zero the tendency array on the GPU

    zero_tend(1)=1

    call oclWriteIntBuffer(zero_tend_buf, zero_tend_sz, zero_tend);

    call runOcl(jikmsz, 0)

! Then compute the new tendencies

    zero_tend(1)=0

    call oclWriteIntBuffer(zero_tend_buf, zero_tend_sz, zero_tend);

    call runOcl(oclGlobalRange, oclLocalRange)

! Read back results from GPU

    call oclReadBuffer(tendency_buf, jik_sz, tendency)

end subroutine advect_scalar_ocl

```

the reason for this smaller improvement is the higher transfer cost on the AMD system. Of course, these indicators ignore the effect of the implementation of the code and the compiler performance, but they give an indication of what is achievable in terms of the hardware capability. What this means is that, if one achieves a higher speed-up than these figures, either the application is not computation dominated, or the coding is sub-optimal, or the compilation is sub-optimal. If the code is memory bandwidth dominated, we see that the achievable speed-up is about $4\times$.

The parallelization of the advection kernel using OpenCL is discussed above, the conclusion is that it is definitely feasible and can result in very good performance. However, we want to focus on the findings from the performance evaluation. In our opinion, the most important finding is that, in order to achieve the best possible performance on either a multicore CPU or a GPU, it is necessary to considerably rewrite the code for data-parallel execution. The other key finding

is that the current system architecture is problematic for GPU acceleration, because the PCIe bus performance constitutes a huge bottleneck. However, it is still possible to achieve good performance provided that a substantial part of the model code is implemented on the GPU.

On the other hand, one has to ask the question if it is at all worthwhile to offload the code to the GPU. To evaluate this question we used both OpenCL and OpenMP to parallelize our kernel code on the multicore CPU. As argued above, there is in fact theoretically almost no difference in performance between the Intel Xeon E5-2640 multicore CPU and the Nvidia Tesla C2070 GPU. In practice, the GPU performance of the scalar advection kernel is worse because of the high cost of moving the data. If we remove this cost – effectively simulating the case of running a fully integrated model on the GPU – we see that the GPU performance is much better than the original (single-threaded) code.

The difference in performance between the GPU and the Intel CPU is a result of a combination of factors: the CPU code is vectorized but single-threaded; the GPU code is parallelised over multiple compute units but the threads within a single compute units can't deliver the theoretical level of parallelism because the memory accesses are not entirely coalesced.

We must introduce another key factor in performance comparisons, often overlooked: the influence of the compiler. For the original Fortran code, we used both the GNU compiler and the commercial PGI compiler, and we found that the code compiled with the latter runs much faster than with the former. The reason is that the PGI compiler makes full use of the Xeon's 256-bit AVX vector instructions, while gcc doesn't. Unfortunately, we did not have a license for the C/C++ version of the PGI compiler and as a result we could not directly evaluate our C++ OpenMP code performance with this compiler. However, it is reasonable to assume that the PGI compiler would produce the same speed-ups for the C++ code as for the Fortran code. As seen from the OpenMP benchmarks (Figure 4), the multi-threaded CPU version could run almost six times faster than the single-threaded version, so compared to that the GPU would be about two times faster. It is unlikely that the CPU would be actually six times faster because due to the vectorization, effectively there

AN ANALYSIS OF THE FEASIBILITY AND BENEFITS OF GPU/MULTICORE ACCELERATION OF WRF31

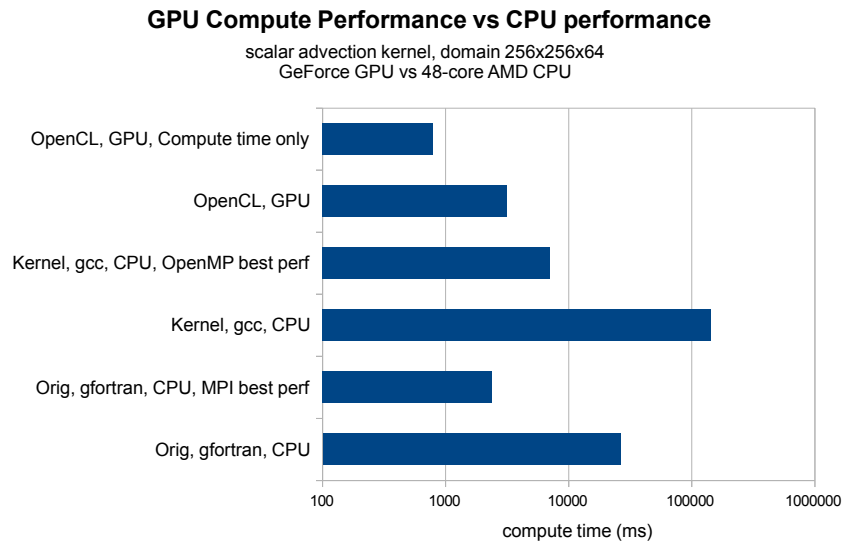


Figure 13. Advection kernel GPU compute performance vs original kernel parallelized using MPI

will be no benefit from hyperthreading, as demonstrated by Saini in [19], so that the expected figure is closer to four times, in other words the GPU would be 3 times faster.

From Table I is is clear that the CPIs for the Intel and AMD 48-core CPUs and the GPU are very similar; furthermore, [20], Langkamp found very little difference in the WRF MPI performance when compiled using gfortran and the PGI compiler (on a AMD Opteron 2384 system). So, using the 48-core AMD system as our reference, we can compare our kernel's compute performance with the original kernel parallelised using MPI. The results are summarised in Figure 13. We see that the OpenCL kernel's compute performance on the GPU is about three times faster than the original kernel parallelized using MPI on the 48-core AMD CPU.

Thus, we can conclude that the performance of the OpenCL kernel on the GPU would be at least three times faster than the parallelized kernel on the CPU, for GPUs and CPUs with comparable CPIs.

Note that if we had only used the GNU compiler on the Intel platform, we would have reported a 20× speed-up; and without compiler optimizations, this figure would be even higher. This example

illustrates the values of the estimates based on the CPIs, as well as illustrating the differences caused by different coding styles and compilers.

A final point concerns the optimization goal: should the code be optimized for speed or for power? For the individual user, the aim is either to reduce the run time of the simulation or increase the accuracy. The limitation for an individual user is usually the cost of purchasing the system, rather than the operating cost. Considering that the low cost of a GeForce GPU, it might be more cost-effective to buy a GPU rather than an additional multicore CPU system.

However, for large high-performance computing centers, the aim is to minimize the energy consumption of the system, because electricity bills are the dominant component in the total cost of ownership. To save energy, one must consider both the power consumption and the speed of execution. For example, if a GPU has the same power consumption as its host CPU, then using it will result in a net energy savings only if the speed-up is greater than a factor of two. Therefore, arguably, the key indicator for assessing hardware acceleration should be the increase in performance-per-Watt. Here, the PCIe-hosted GPU is at a disadvantage because it can't work without its host, and even in idle mode the power consumption of a large multicore CPU is considerable. Hosting the GPU on a low-power ARM or Atom based system is a possible option to alleviate this issue. A GPU-CPU hybrid such as the AMD Fusion or the low-power Nvidia Tegra 4 could potentially be an even better choice.

8. CONCLUSIONS

Based on our work we can conclude that GPU acceleration of NWP codes such as WRF is both feasible and worthwhile, but that a number of important issues remain to be addressed.

A very important conclusion is that rewriting the code as OpenCL-style data-parallel kernels can already result in significant speed-up of the code on a multicore CPU system using either OpenCL or OpenMP, i.e. without using a GPU. Consequently, this is an essential step. However, in particular for WRF this requires a major rewrite of the dynamics and physics code.

AN ANALYSIS OF THE FEASIBILITY AND BENEFITS OF GPU/MULTICORE ACCELERATION OF WRF33

Another important finding is that the current PCIe-based CPU-GPU system architecture is sub-optimal for NWP acceleration because of the huge bottleneck of the data transfers over the PCIe bus. On the one hand, this means that a considerable part of the code must be executed on the GPU to amortize this cost. On the other hand, it means that the new CPU-GPU hybrid chips could be very promising for NWP acceleration.

Our final assessment is that even with the current system architectures, accelerating WRF with a factor of up to five times is definitely an achievable goal.

It is important to note that our findings are more generally applicable to multi-physics fluid dynamics codes: in many fluid dynamics codes the numerical schemes of the advection terms are based on finite differences between neighboring cells, similar to the WRF code. For fluid systems including multi-physics processes, there are many calls to these advection routines. This class of numerical codes will benefit from hardware acceleration.

Accelerating multi-physics fluid dynamics codes including NWP codes is critically important for forecasting applications in atmospheric and environmental issues. Forecasting of extreme weather events, early warning of environmental pollution, and emergency response to the dispersion of hazardous materials all requires fast and accurate computations of multi-physics atmospheric motion. For example, a numerical forecasting of micro-scale atmospheric motions in urban areas and/or over complex topography should be benefitted from computational accelerations, because it requires a coupling approach merging NWP and CFD codes [21] or very high resolutions to accurately represent complex topography [22, 23]. Furthermore, the computational accelerations would be advantageous in climate prediction simulations with high-resolution global- and regional-scale atmospheric model for better representing tropical cyclones and heavy rainfall systems [24, 25].

REFERENCES

1. Shainer G, Liu T, Michalakes J, Liberman J, Layton J, Celebioglu O, Schultz SA, Mora J, Cownie D. Weather research and forecast (wrf) model performance and profiling analysis on advanced multi-core hpc clusters. *The*

- 10th LCI International Conference on High-Performance Clustered Computing, Boulder, CO 2009*; .
2. Morton D, Nudson O, Stephenson C. Benchmarking and evaluation of the weather research and forecasting (wrf) model on the cray xt5. *Cray User Group Proceedings, Atlanta, GA 2009*; :04–07.
 3. Michalakes J, Dudhia J, Gill D, Henderson T, Klemp J, Skamarock W, Wang W. The weather research and forecast model: software architecture and performance. *Proceedings of the 11th ECMWF Workshop on the Use of High Performance Computing In Meteorology*, vol. 25, World Scientific, 2004; 29.
 4. World Economic Forum. Global Risks 2011 Sixth Edition – An initiative of the Risk Response Network Jan 2011.
 5. Stern N. *The economics of climate change: the Stern review*. Cambridge Univ Pr, 2007.
 6. Rosenzweig C, Iglesias A, Yang X, Epstein P, Chivian E. Climate change and us agriculture: The impacts of warming and extreme weather events on productivity, plant diseases, and pests. *Center for Health and the Global Environment, Harvard Medical School, Boston, MA, USA 2000*; .
 7. Luber G, McGeehin M. Climate change and extreme heat events. *American Journal of Preventive Medicine* 2008; **35**(5):429–435.
 8. Knutti R, Furrer R, Tebaldi C, Cermak J, Meehl G. Challenges in combining projections from multiple climate models. *Journal of Climate* 2010; **23**(10):2739–2758.
 9. Michalakes J, Chen S, Dudhia J, Hart L, Klemp J, Middlecoff J, Skamarock W. Development of a next generation regional weather research and forecast model. *Developments in Teracomputing: Proceedings of the Ninth ECMWF Workshop on the Use of High Performance Computing in Meteorology*, vol. 1, World Scientific, 2001; 269–276.
 10. Skamarock WC, Klemp JB, Dudhia J, Gill DO, Barker DM, Duda MG, Huang XY, Wang W, Powers JG. A description of the Advanced Research WRF Version 3. *Technical Report*, NCAR/TN-475+STR 2008.
 11. Skamarock WC, Klemp JB. A time-split nonhydrostatic atmospheric model for weather research and forecasting applications. *Journal of Computational Physics* 2008; **227**:3465–3485.
 12. Michalakes J, Vachharajani M. GPU acceleration of numerical weather prediction. *Parallel and Distributed Processing, 2008. IPDPS 2008. IEEE International Symposium on*, IEEE, 2008; 1–7.
 13. Govett M, Middlecoff J, Henderson T, Rosinski J, Madden P. Successes and Challenges Porting Weather and Climate Models to GPUs. *AGU Fall Meeting Abstracts*, vol. 1, 2011; 02.
 14. Fuhrer O, Gysi T, Lapillonne X, Osuna C, Cumming B, Sawyer W, Messme P, , Schroeder T, Schulthess TC. GPU Consideration for Next Generation Weather and Climate Simulations. *Technical Report*, CSCS Swiss National Supercomputing Centre 2012.
 15. Shimokawabe T, Aoki T, Muroi C, Ishida J, Kawano K, Endo T, Nukada A, Maruyama N, Matsuoka S. An 80-fold speedup, 15.0 TFlops full GPU acceleration of non-hydrostatic weather model ASUCA production code. *High Performance Computing, Networking, Storage and Analysis (SC), 2010 International Conference for*, IEEE, 2010; 1–11.
 16. Munshi A, *et al.*. The OpenCL Specification. *Technical Report*, Khronos OpenCL Working Group 2009.
 17. Hong SY, Lim JO. The wrf single-moment 6-class microphysics scheme (wsm6). *Journal of the Korean Meteorological Society* 2006; **42**:129–151.

AN ANALYSIS OF THE FEASIBILITY AND BENEFITS OF GPU/MULTICORE ACCELERATION OF WRF 35

18. Porter A, Ashworth M, Gadian A, Burton R, Connolly P, Bane M. WRF Code Optimisation for Meso-Scale Process Studies (WOMPS) dCSE Project Report. *Technical Report*, STFC 2012.
19. Saini S, Jin H, Hood R, Barker D, Mehrotra P, Biswas R. The impact of hyper-threading on processor resource utilization in production applications. *High Performance Computing (HiPC), 2011 18th International Conference on*, IEEE, 2011; 1–10.
20. Langkamp T, Böhner J. Influence of the compiler on multi-cpu performance of wrfv3. *Geoscientific Model Development* 2011; **4**(3):611–623.
21. Nakayama H, Takemi T, Nagai H. Large-eddy simulation of urban boundary-layer flows by generating turbulent inflows from mesoscale meteorological simulations. *Atmospheric Science Letters* 2012; **13**(3):180–186.
22. Takemi T. High-resolution meteorological simulations of local-scale wind fields over complex terrain: A case study for the eastern area of fukushima in march 2011. *Theoretical and Applied Mechanics Japan* 2013; **61**:3–10.
23. Oku Y, Takemi T, Ishikawa H, Kanada S, Nakano M. Representation of extreme weather during a typhoon landfall in regional meteorological simulations: A model intercomparison study for typhoon songda (2004). *Hydrologic Research Letters* 2010; **4**:1–5.
24. Mizuta R, Yoshimura H, Murakami H, Matsueda M, Endo H, Ose T, Kamiguchi K, Hosaka M, Sugi M, Yukimoto S, *et al.*. Climate simulations using mri-agcm3.2 with 20-km grid. *Journal of the Meteorological Society of Japan* 2012; **90A**:233–258.
25. Kanada S, Nakano M, Kato T. Projections of future changes in precipitation and the vertical structure of the frontal zone during the baiu season in the vicinity of japan using a 5-km-mesh regional climate model. *Journal of the Meteorological Society of Japan* 2012; **90A**:65–86.